

haberman dataset EDA plotting

```
In [1]: import warnings
warnings.filterwarnings("ignore")

In [20]: #haberman dataset EDA plotting
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
#load haberman.csv into a pandas dataframe.
hman = pd.read_csv("haberman.csv")

In [21]: #tells the shape of given csv file
print(hman.shape)

(306, 4)

In [22]: #columns present in file
print(hman.columns)

Index(['Age', 'Op_Year', 'axil_nodes_det', 'Surv_status'], dtype='object')

In [5]: #no. of patients alive after 5 year of cancer treatment and patients who cant survive
hman["Surv_status"].value_counts()

Out[5]: 1    225
        2     81
        Name: Surv_status, dtype: int64

In [24]: hman["Surv_status"][hman["Surv_status"]==1]="Yes"
hman["Surv_status"][hman["Surv_status"]==2]="No"
```

observation

- 1. Here "surv\_counts" shows the survival of patients after 5 year of cancer treatment.
- 2. "1" shows that patient survive 5 year or more .
- 3. "2" shows that patient can not survive 5 year.
- 4. according to this data there are 225 people survives and 81 can not survive.

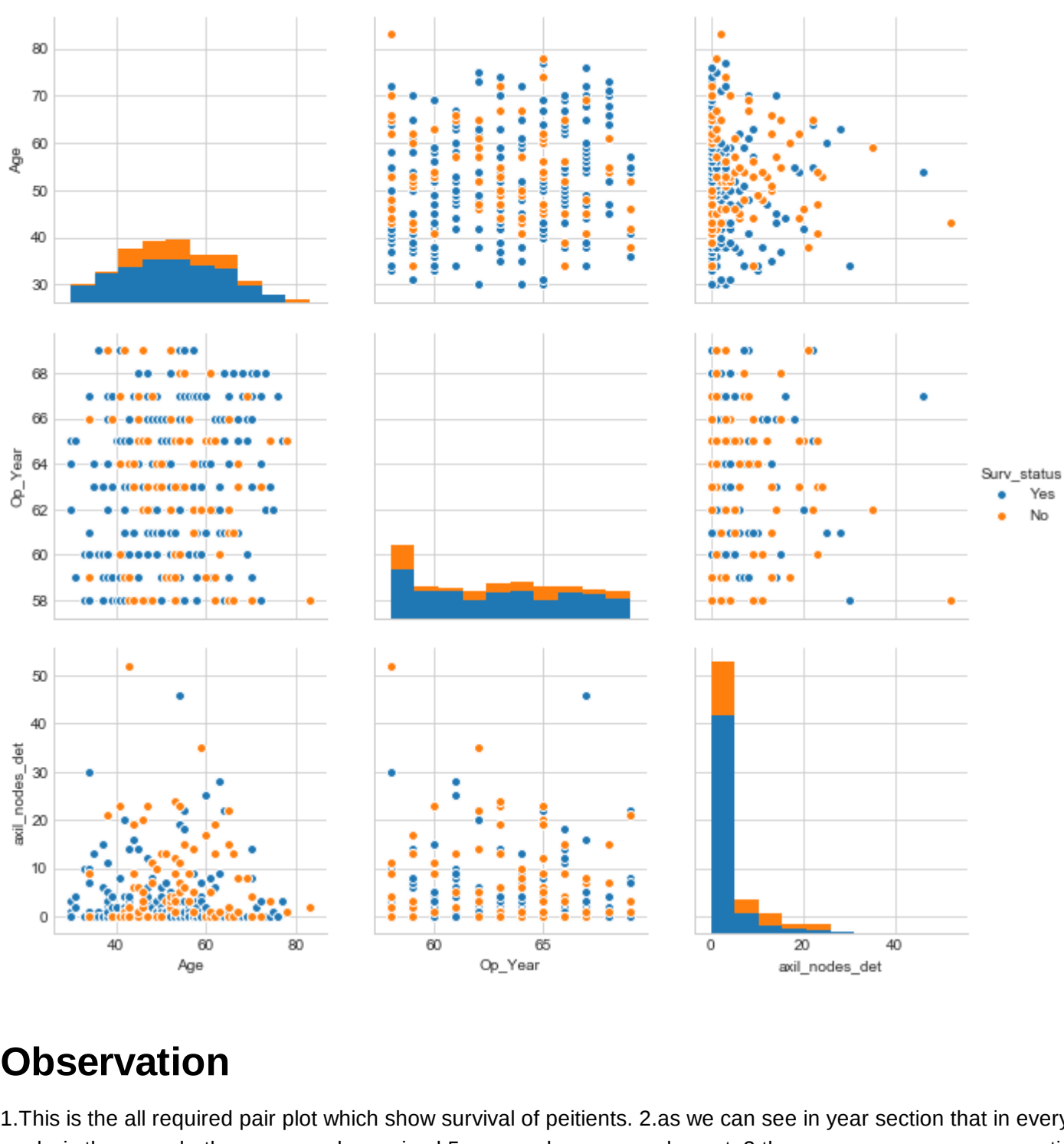
2-D scatter plot

```
In [11]: #2-D scatter plot of survey status of patients w.r.t. axillary nodes detected
hman.plot(kind="scatter",x="axil_nodes_det",y="Surv_status")
plt.show()
```

```
In [12]: #2-D scatter plot of survey status of patients w.r.t. axillary nodes detected in different hue
import seaborn as sns
sns.set_style("whitegrid")
sns.FacetGrid(hman,hue="Surv_status",size=4)\
    .map(plt.scatter,"axil_nodes_det","Surv_status")\
    .add_legend()
plt.show()
```

Pair plot

```
In [26]: #pair plot of the give data of hyberman
plt.close()
sns.set_style("whitegrid")
sns.pairplot(hman,hue="Surv_status",size=3)
plt.show()
```



Observation

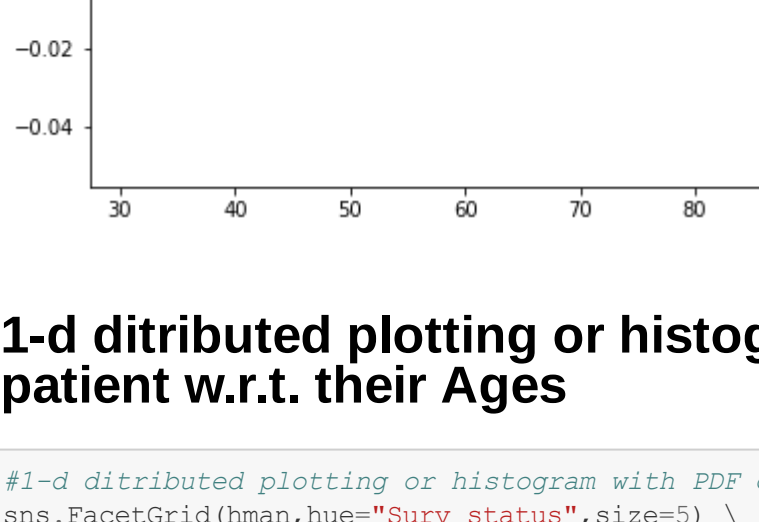
1.This is the all required pair plot which show survival of petients. 2.as we can see in year section that in every year during analysis there are both some people survived 5 year and some people cant. 3.there are some more assumptions we can make after analysing these plottings

1-D plotting

```
In [8]: #1-D scatter plot of Age of patients with respect to their survey status
import numpy as np
hman_surv = hman.loc[hman["Surv_status"] == 1]
hman_nsurv = hman.loc[hman["Surv_status"] == 2]

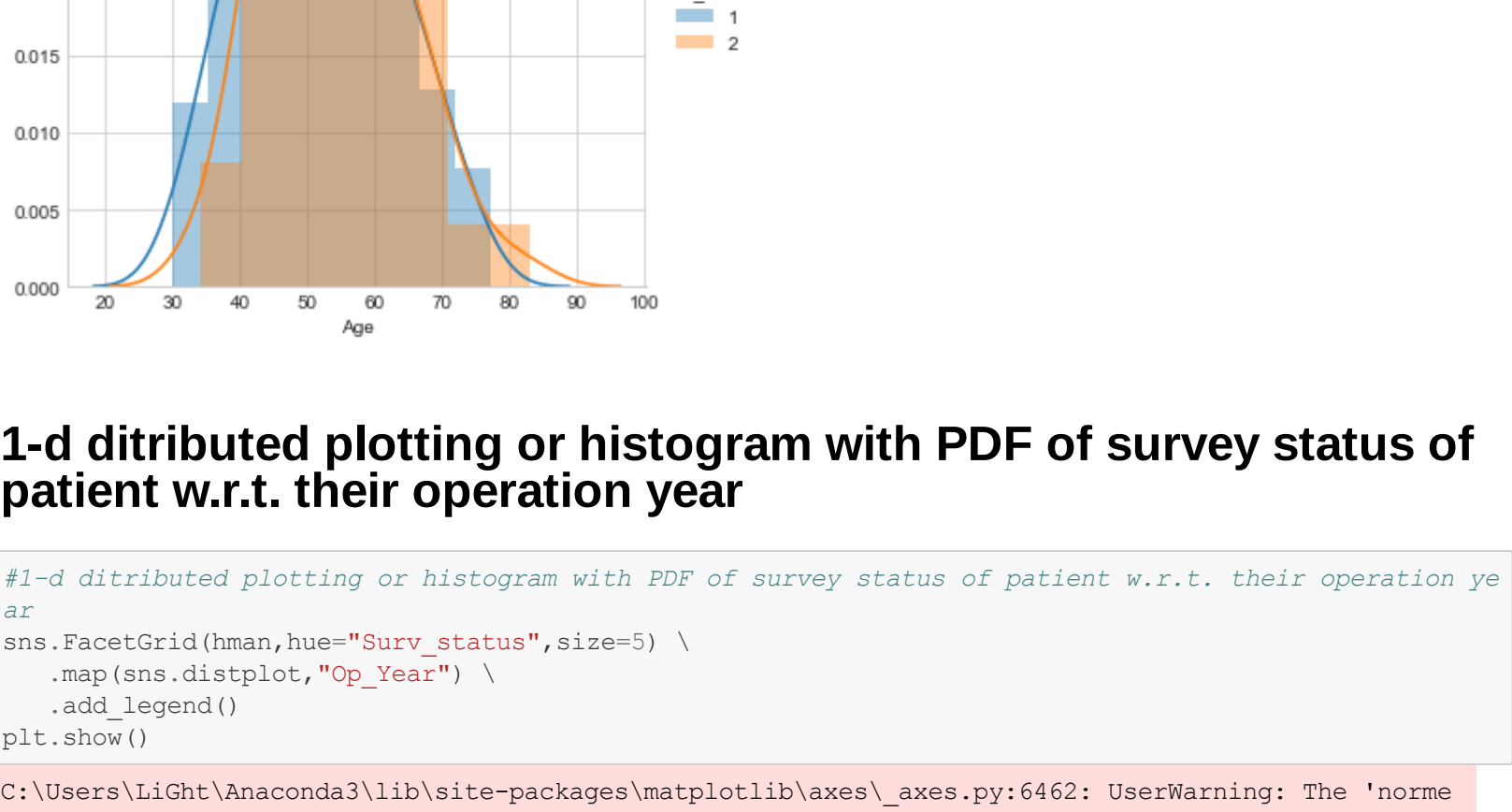
plt.plot(hman_surv["Age"],np.zeros_like(hman_surv['Age']),"o")
plt.plot(hman_nsurv["Age"],np.zeros_like(hman_nsurv['Age']),"o")

plt.show()
```



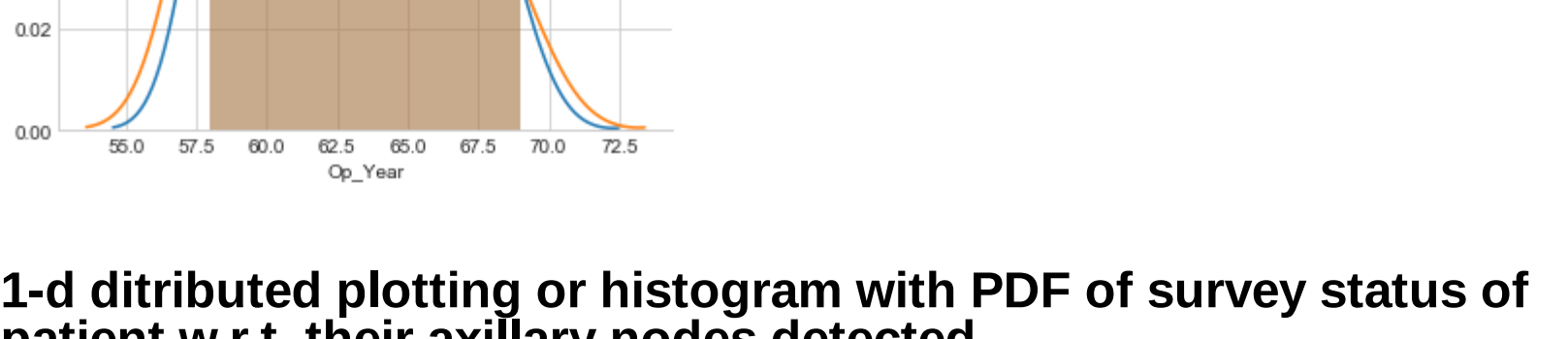
1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their Ages

```
In [21]: #1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their Ages
sns.FacetGrid(hman,hue="Surv_status",size=5) \
    .map(sns.distplot,"Age") \
    .add_legend()
plt.show()
```



1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their operation year

```
In [25]: #1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their operation ye
ar
sns.FacetGrid(hman,hue="Surv_status",size=5) \
    .map(sns.distplot,"Op_Year") \
    .add_legend()
plt.show()
```

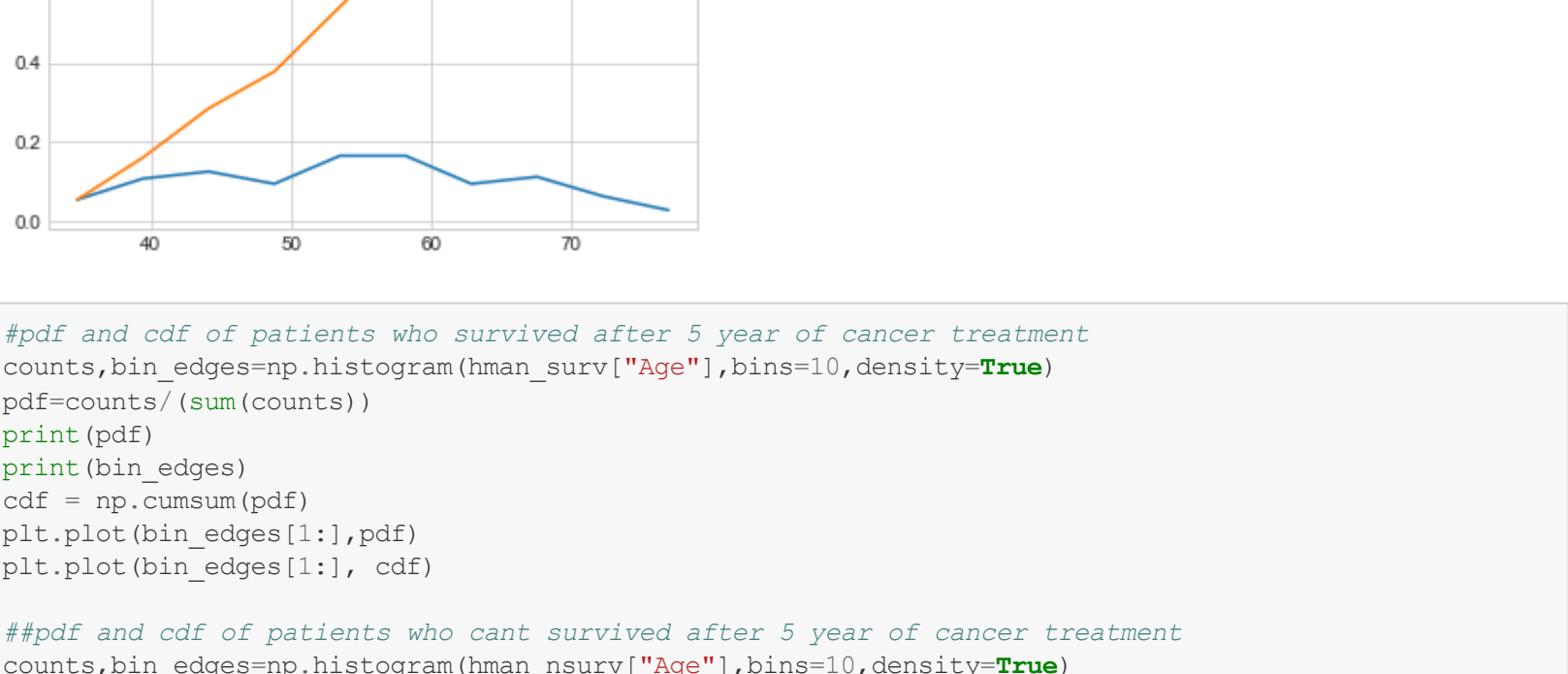


1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their axillary nodes detected

```
In [ ]: # 1-d ditributed plotting or histogram with PDF of survey status of patient w.r.t. their axillary no
des detected
sns.FacetGrid(hman,hue="Surv_status",size=5) \
    .map(sns.distplot,"axil_nodes_det") \
    .add_legend()
plt.show()
```

pdf and cdf of patients who survived after 5 year of cancer treatment

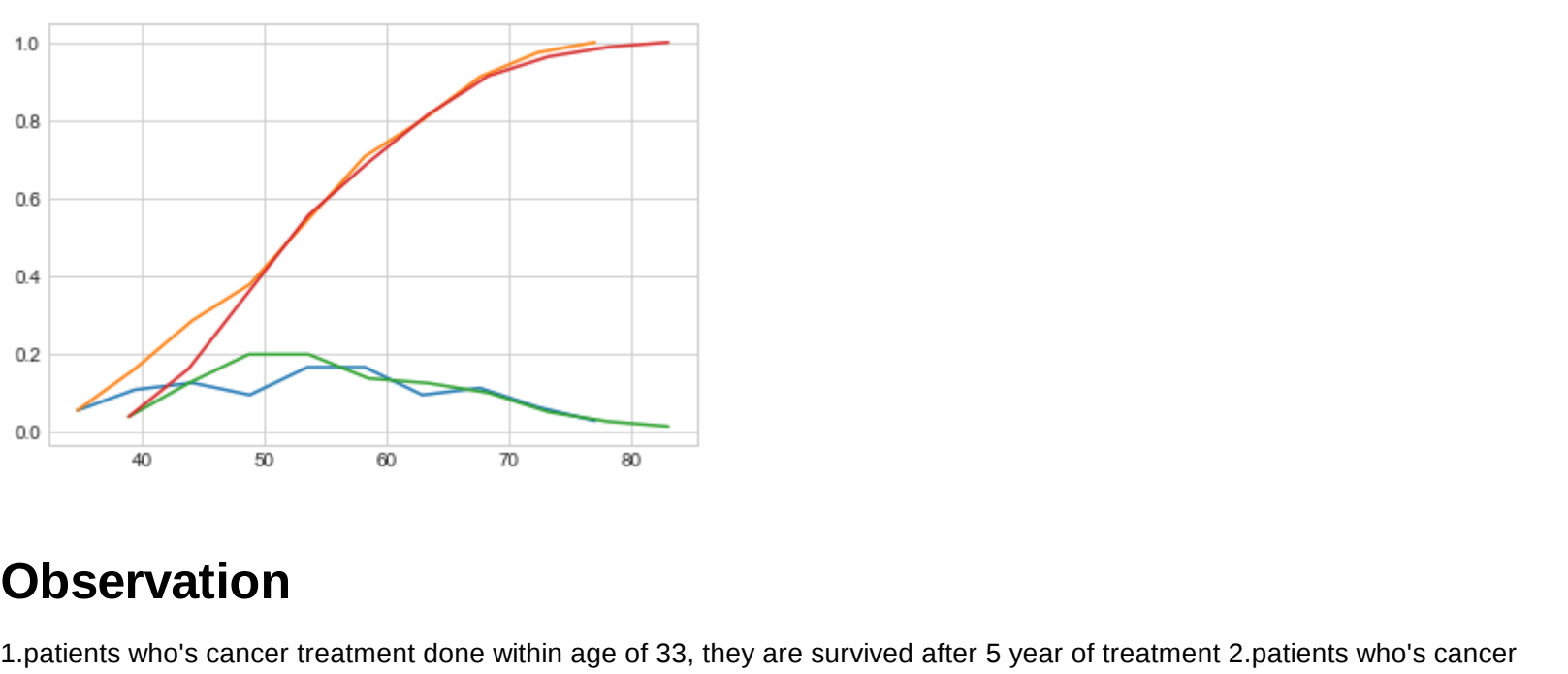
```
In [23]: # PDF and CDF of given data
#pdf and cdf of patients who survived after 5 year of cancer treatment
counts,bin_edges=np.histogram(hman_surv["Age"],bins=10,density=True)
pdf=counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:],cdf)
plt.show()
```



```
In [25]: #pdf and cdf of patients who survived after 5 year of cancer treatment
counts,bin_edges=np.histogram(hman_surv["Age"],bins=10,density=True)
pdf=counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:],cdf)

#pdf and cdf of patients who cant survived after 5 year of cancer treatment
counts,bin_edges=np.histogram(hman_nsurv["Age"],bins=10,density=True)
pdf=counts/(sum(counts))
print(pdf)
print(bin_edges)
cdf = np.cumsum(pdf)
plt.plot(bin_edges[1:],pdf)
plt.plot(bin_edges[1:],cdf)

plt.show()
```



Observation

1.patients who's cancer treatment done within age of 33, they are survived after 5 year of treatment 2.patients who's cancer treatment done after age of 78, they can't survived after 5 year of treatment

Mean, variance and standard deviation

```
In [33]: # mean of given data
print("Means:")
print(np.mean(hman_surv["Age"]))
print(np.mean(hman_nsurv["Age"]))
# standard deviation of given data
print("Unstandard deviation")
print(np.std(hman_surv["Age"]))
print(np.std(hman_nsurv["Age"]))

Means:
52.01777777777778
53.67901234567901

standard deviation:
10.98765547510051
10.10418219303131
```

Median, Percentile, Quantile, IQR, MAD

```
In [41]: # median of given data
print("medians")
print(np.median(hman_surv["Age"]))
print(np.median(hman_nsurv["Age"]))
# median of given data
print("\n 90th Percentile:")
print(np.percentile(hman_surv["Age"],90))
print(np.percentile(hman_nsurv["Age"],90))
# quantiles of given data
print("\nQuantiles:")
print(np.percentile(hman_surv["Age"],np.arange(0,100,25)))
print(np.percentile(hman_nsurv["Age"],np.arange(0,100,25)))
# MAD of given data
from statsmodels import robust
print("\nMAD:")
print(robust.mad(hman_surv["Age"]))
print(robust.mad(hman_nsurv["Age"]))

medians:
52.0
53.0

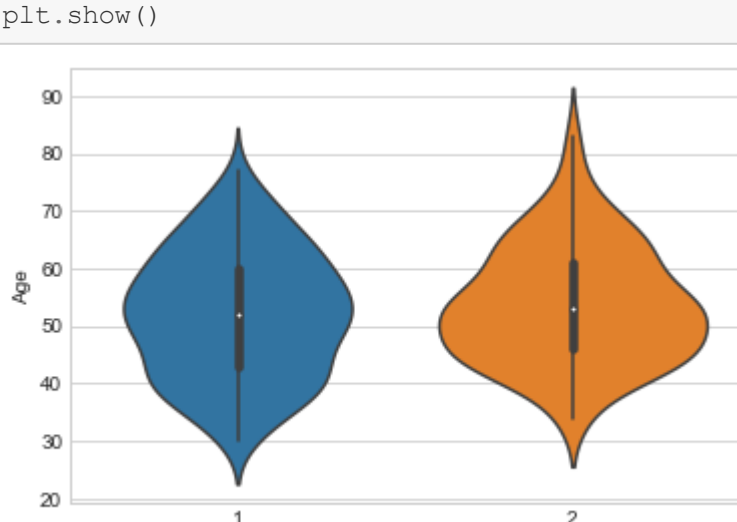
90th Percentile:
67.0
67.0

quantiles:
[30. 43. 52. 60.]
[34. 46. 53. 61.]

MAD:
13.343419966550417
11.860817748044816
```

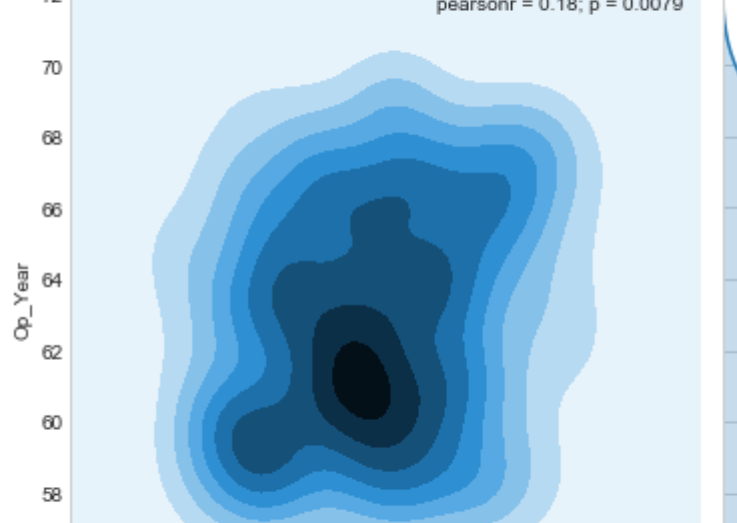
Box plot and Whiskers

```
In [45]: #box plot and whiskers to show survey
sns.boxplot(x="Surv_status",y="Age",data=hman)
plt.show()
```



Violin plots

```
In [46]: #violin plot of given haberman data
sns.violinplot(x="Surv_status",y="Age",data=hman,size=10)
plt.show()
```



2-D contour plot

```
In [47]: #2D Density plot, contors-plot
sns.jointplot(x="Age", y="Op_Year", data=hman_surv, kind="kde");
plt.show();
```



Final Observation